## Nonlinear optimization

without constraints

Michel Bierlaire

Introduction to optimization and operations research



## Nonlinear optimization

#### Problem definition

$$\min_{x \in \mathbb{R}^n} f(x)$$

where f is twice differentiable and bounded from below:

$$\exists M \in \mathbb{R} \text{ such that } f(x) \geq M, \ \forall x \in \mathbb{R}^n.$$

### Fermat's theorem

#### Theorem 5.1

- $\triangleright$   $x^*$  is a local minimum of  $f: \mathbb{R}^n \to \mathbb{R}$ .
- ▶ If f is differentiable around  $x^*$ , then

$$\nabla f(x^*) = 0.$$

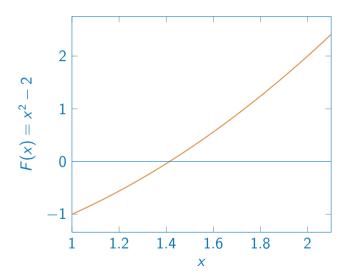
▶ If f is twice differentiable around  $x^*$ , then

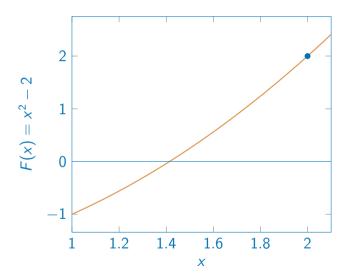
$$\nabla^2 f(x^*) \ge 0$$
 [positive semidefinite].

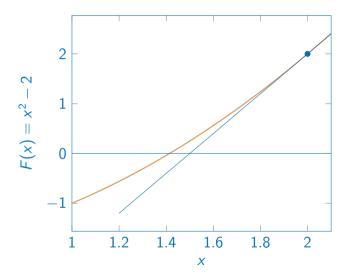
## Solving one equation with one variable

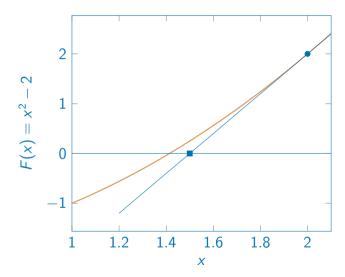
#### Motivation

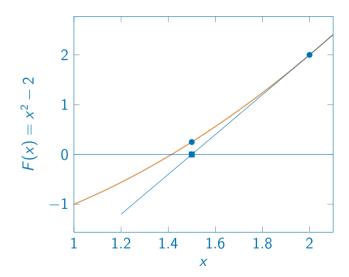
- Optimization algorithms rely on solving systems of equations given by optimality conditions.
- Systems of linear equations can be solved by Gaussian elimination.
- Systems of nonlinear equations can be solved using Newton's method.
- ▶ We remind this method on the simple case of one equation and one variable.

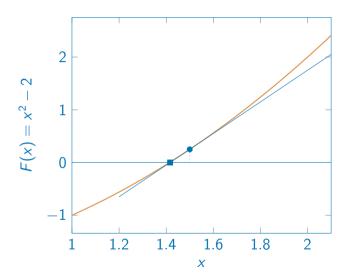












### **Iterations**

k	$x_k$	$F(x_k)$	$F'(x_k)$
0	+2.00000000E+00	+2.00000000E+00	+4.0000000E+00
1	+1.50000000E+00	+2.50000000E-01	+3.0000000E+00
2	+1.41666667E+00	+6.9444444E-03	+2.8333333E+00
3	+1.41421569E+00	+6.00730488E-06	+2.82843137E+00
4	+1.41421356E+00	+4.51061410E-12	+2.82842712E+00
5	+1.41421356E+00	+4.44089210E-16	+2.82842712E+00

# Speed of convergence



- ► Fast method.
- ► The precision doubles at each iteration.

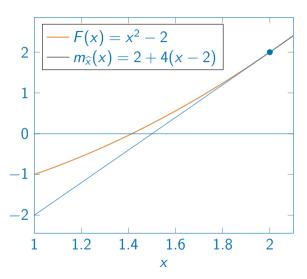
### Linear model

### Taylor's theorem

$$F: \mathbb{R} \to \mathbb{R}$$
 differentiable  $m_{\widehat{x}}(x) = F(\widehat{x}) + (x - \widehat{x})F'(\widehat{x})$ 

$$F(x) = x^2 - 2F'(x) = 2x$$

$$\hat{x} = 2 : m_{\hat{x}}(x) = 2 + 4(x - 2)$$



## Algorithm

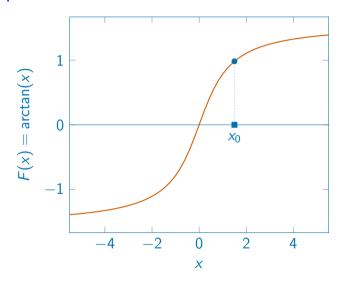
At each iteration, find  $x_{k+1}$  that solves

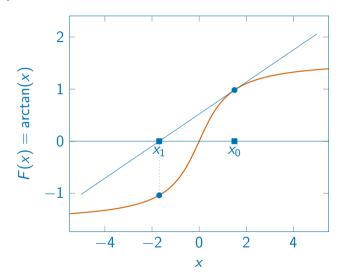
$$F(x_k) + (x - x_k)F'(x_k) = 0$$

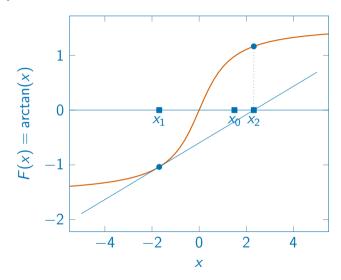
$$(x_{k+1} - x_k)F'(x_k) = -F(x_k)$$

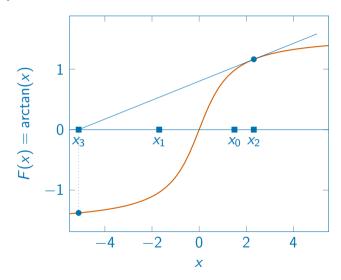
$$x_{k+1} - x_k = -F(x_k)/F'(x_k)$$

$$x_{k+1} = x_k - F(x_k)/F'(x_k)$$









## Convergence



- ► The method is fast...
- but does not always converge.
- Several conditions have to be met.

### Convergence

#### Motivation

- ▶ We have seen that Newton's method, when it converges, does so fast.
- But it does not always converge.
- ▶ We see now in what circumstances it converges, and quantify its speed.
- ▶ Then, we generalize the method for systems of equations.

#### Conditions

F is not too non linear

F' is Lipschitz continuous, with Lipschitz constant M.

 $\exists M > 0$  such that  $\forall x, y, |F'(x) - F'(y)| \leq M|x - y|$ .

#### **Conditions**

#### F is not too non linear

F' is Lipschitz continuous, with Lipschitz constant M.

$$\exists M > 0$$
 such that  $\forall x, y, |F'(x) - F'(y)| \leq M|x - y|$ .

F' is not too close to 0

$$\exists \rho > 0 \text{ such that } |F'(x)| \geq \rho, \forall x.$$

#### **Conditions**

#### F is not too non linear

F' is Lipschitz continuous, with Lipschitz constant M.

$$\exists M > 0$$
 such that  $\forall x, y, |F'(x) - F'(y)| \leq M|x - y|$ .

F' is not too close to 0

$$\exists \rho > 0$$
 such that  $|F'(x)| \geq \rho, \forall x$ .

 $x_0$  is not too far from the solution

$$\exists \eta > 0 \text{ such that } |x_0 - x^*| < \eta.$$

## Convergence

#### Theorem 7.7

Consider the sequence of iterates:

$$x_{k+1} = x_k - \frac{F(x_k)}{F'(x_k)}, \qquad k = 0, 1, \ldots,$$

- ▶ it is well defined,
- $ightharpoonup \lim_{k\to\infty} x_k = x^*$ ,
- ▶ it converges *q*-quadratically:

$$|x_{k+1}-x^*| \leq \frac{M}{2\rho} |x_k-x^*|^2$$
.

### Extension to *n* variables

#### Problem statement

$$F: \mathbb{R}^n \to \mathbb{R}^n$$
 differentiable

Find  $x \in \mathbb{R}^n$  such that

$$F(x)=0.$$

### Reminder

$$F: \mathbb{R}^n \to \mathbb{R}^n$$

Gradient matrix

$$\nabla F(x) = \begin{pmatrix} \begin{vmatrix} & & & \\ \nabla F_1(x) & \cdots & \nabla F_n(x) \\ & & \end{vmatrix} \end{pmatrix} \qquad J(x) = \nabla F(x)^T$$

$$= \begin{pmatrix} \frac{\partial F_1}{\partial x_1} & \frac{\partial F_2}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_1} \\ \vdots & \vdots & \vdots & \vdots \\ \frac{\partial F_1}{\partial x_n} & \frac{\partial F_2}{\partial x_n} & \cdots & \frac{\partial F_n}{\partial x_n} \end{pmatrix} \qquad = \begin{pmatrix} ---- & \nabla F_1(x)^T & ---- \\ & \vdots & \vdots \\ ---- & \nabla F_n(x)^T & ---- \end{pmatrix}.$$

$$= \begin{pmatrix} \cdots & \nabla F_1(x)^T & \cdots \\ & \vdots & \\ \cdots & \nabla F_n(x)^T & \cdots \end{pmatrix}.$$

### Linear model

$$m_{\widehat{x}}(x) = F(\widehat{x}) + \nabla F(\widehat{x})^{T}(x - \widehat{x}) = F(\widehat{x}) + J(\widehat{x})(x - \widehat{x})$$

### **Iterations**

$$x_{k+1} = x_k + d_k$$

where

$$J(x_k)d_k=-F(x_k).$$

## Solving the necessary optimality conditions

#### Motivation

- ► The necessary optimality condition stipulates that, at a local optimum, the gradient is zero.
- ▶ The first algorithm consists in solving the system of equations  $\nabla f(x) = 0$  in order to find a stationary point...
- ... using Newton's method.

### Newton's method

### Equations

### Problem

$$F(x) = 0$$

### Algorithm

$$x_{k+1} = x_k + d_k$$

where  $d_k$  is solution of

$$J(x_k)d_k=-F(x_k).$$

### Optimization

### **Problem**

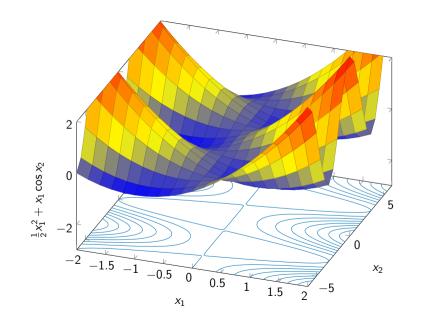
$$\nabla f(x) = 0$$

### Algorithm

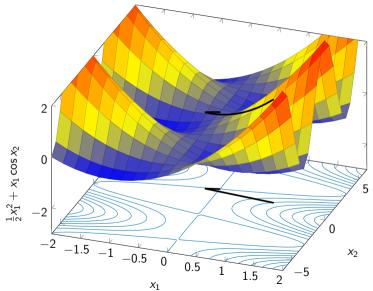
$$x_{k+1} = x_k + d_k$$

where  $d_k$  is solution of

$$\nabla^2 f(x_k) d_k = -\nabla f(x_k).$$



Example: fast convergence



## **Iterations**

k	$x_k$	$\nabla f(x_k)$	$\ \nabla f(x_k)\ $	$f(x_k)$
0	1.0000000e+00	1.54030230e+00	1.75516512e+00	1.04030231e+00
	1.00000000e+00	-8.41470984e-01		
1	-2.33845128e-01	-2.87077027e-02	2.30665381e-01	7.53121618e-02
	1.36419220e+00	2.28871986e-01		
2	1.08143752e-02	-3.22524807e-03	1.12840544e-02	-9.33543838e-05
	1.58483641e+00	-1.08133094e-02		
3	-2.13237666e-06	9.22828706e-07	2.32349801e-06	8.79175320e-12
	1.57079327e+00	2.13237666e-06		
4	1.99044272e-17	8.11347449e-17	8.35406072e-17	1.35248527e-25
	1.57079632e+00	-1.99044272e-17		

## Solution

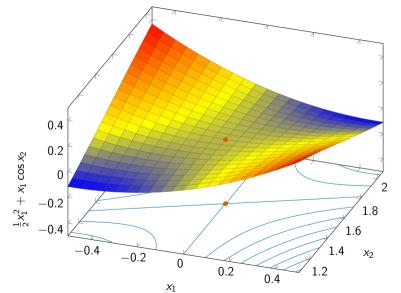
$$x^* = \begin{pmatrix} 0 \\ \pi/2 \end{pmatrix}$$

 $\nabla f(x^*) = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ 

$$abla^2 f(x^*) = \left(egin{array}{cc} 1 & -1 \ -1 & 0 \end{array}
ight)$$

$$\lambda_1 = -0.61803, \ \lambda_2 = 1.6180$$

# Saddle point



#### Geometric interpretation

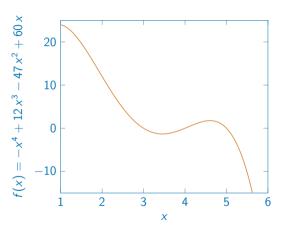
#### Motivation

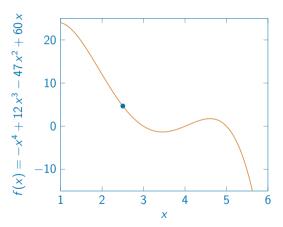
- ▶ In the context of equations, Newton's method uses at each iteration a linear model of the function.
- ▶ In the context of optimization, it uses a linear model of the gradient of the objective function.
- ▶ Equivalently, it uses a quadratic model of the objective function.

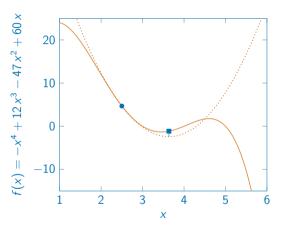
#### Quadratic model

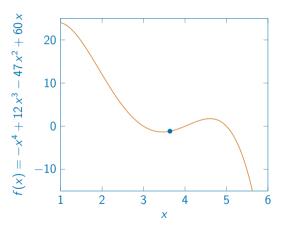
$$f: \mathbb{R}^n \to \mathbb{R}$$
 twice differentiable

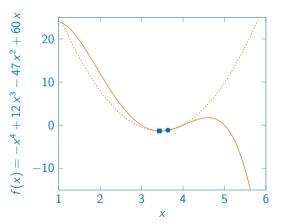
$$m_{\widehat{x}}(x) = f(\widehat{x}) + (x - \widehat{x})^T \nabla f(\widehat{x}) + \frac{1}{2} (x - \widehat{x})^T \nabla^2 f(\widehat{x}) (x - \widehat{x})$$

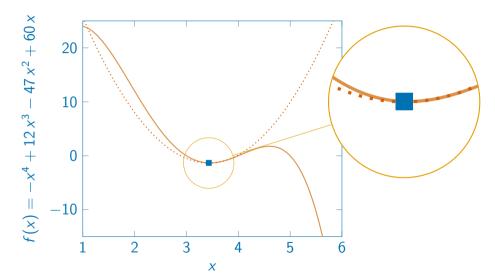


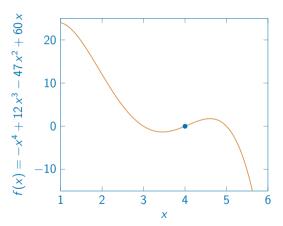


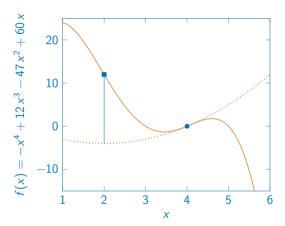


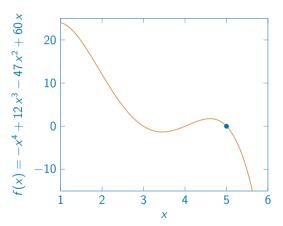


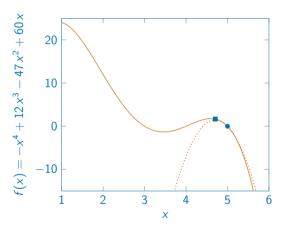












#### Newton's point

- $ightharpoonup f: \mathbb{R}^n \to \mathbb{R}$  twice differentiable
- $ightharpoonup x_k \in \mathbb{R}^n$  such that  $\nabla^2 f(x_k) > 0$

The Newton's point of f at  $x_k$  is

$$x_N = x_k + d_N \,,$$

where  $d_N$  verifies Newton's equations:

$$\nabla^2 f(x_k) d_N = -\nabla f(x_k).$$

### Cauchy's point

- $ightharpoonup f: \mathbb{R}^n \to \mathbb{R}$  twice differentiable
- $ightharpoonup x_k \in \mathbb{R}^n$  such that  $\nabla f(x_k)^T \nabla^2 f(x_k) \nabla f(x_k) > 0$

The Cauchy point of f at  $x_k$  is

$$x_C = x_k - \alpha_C \nabla f(x_k),$$

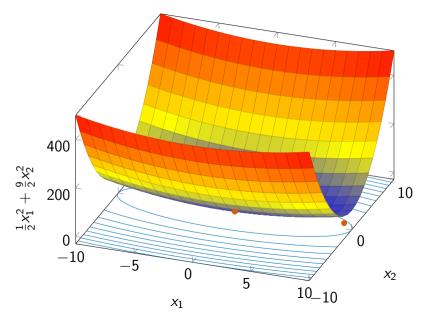
where

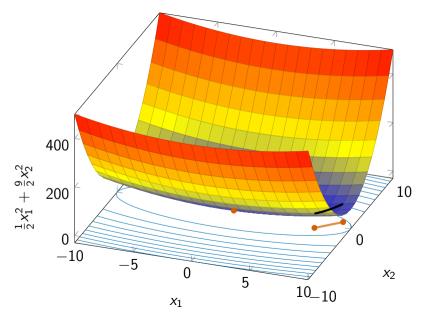
$$\alpha_C = \frac{\nabla f(x_k)^T \nabla f(x_k)}{\nabla f(x_k)^T \nabla^2 f(x_k) \nabla f(x_k)}.$$

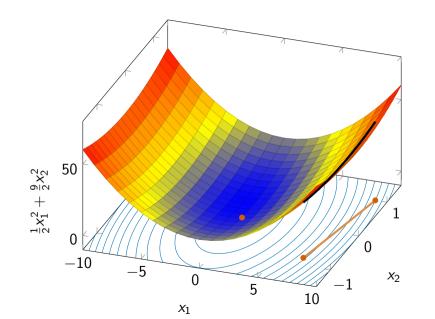
#### Preconditioned steepest descent

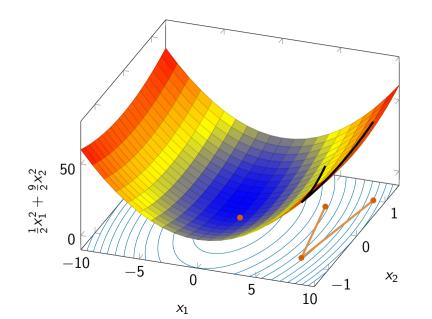
#### Motivation

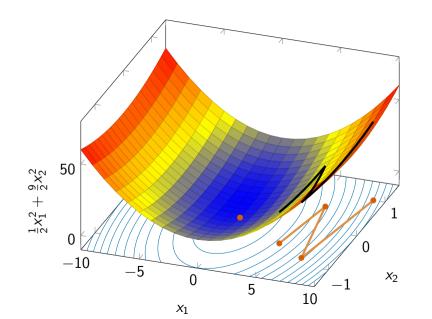
- ▶ A general idea to minimize a function consists in
  - finding a direction along which the function decreases,
  - follow that direction to find a point with a lower value of the objective function.
- This family of methods is called descent methods.
- ► A natural direction to follow is the negative gradient, as it is the steepest descent direction.

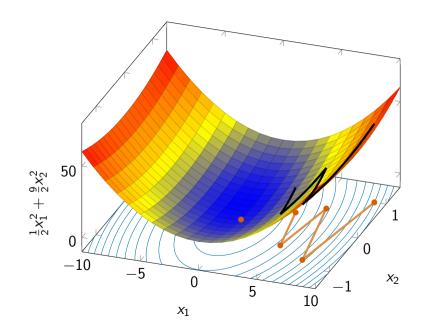


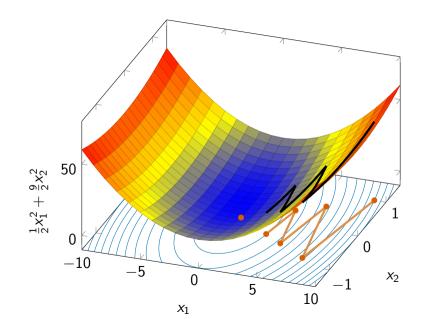












#### **Iterations**

k	$(x_k)_1$	$(x_k)_2$	$\nabla f(x_k)_1$	$\nabla f(x_k)_2$	$\alpha_k$	$f(x_k)$
0	+9.000000E+00	+1.000000E+00	+9.000000E+00	+9.000000E+00	0.2	+4.500000E+01
1	+7.200000E+00	-8.000000E-01	+7.200000E+00	-7.200000E+00	0.2	+2.880000E+01
2	+5.760000E+00	+6.400000E-01	+5.760000E+00	+5.760000E+00	0.2	+1.843200E+01
3	+4.608000E+00	-5.120000E-01	+4.608000E+00	-4.608000E+00	0.2	+1.179648E+01
4	+3.686400E+00	+4.096000E-01	+3.686400E+00	+3.686400E+00	0.2	+7.549747E+00
5	+2.949120E+00	-3.276800E-01	+2.949120E+00	-2.949120E+00	0.2	+4.831838E+00
:						
50	+1.284523E-04	+1.427248E-05	+1.284523E-04	+1.284523E-04	0.2	+9.166662E-09
51	+1.027618E-04	-1.141798E-05	+1.027618E-04	-1.027618E-04	0.2	+5.86664E-09
52	+8.220947E-05	+9.134385E-06	+8.220947E-05	+8.220947E-05	0.2	+3.754665E-09
53	+6.576757E-05	-7.307508E-06	+6.576757E-05	-6.576757E-05	0.2	+2.402985E-09
54	+5.261406E-05	+5.846007E-06	+5.261406E-05	+5.261406E-05	0.2	+1.537911E-09
55	+4.209125E-05	-4.676805E-06	+4.209125E-05	-4.209125E-05	0.2	+9.842628E-10

#### Preconditioning

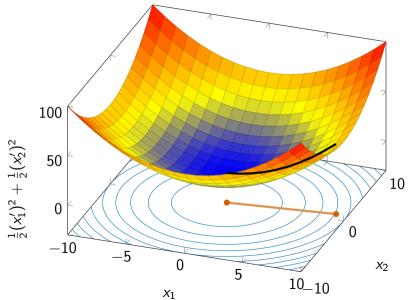
#### Original problem

$$\min_{x \in \mathbb{R}^2} f(x) = \frac{1}{2} x_1^2 + \frac{9}{2} x_2^2$$

#### Change of variables

$$egin{aligned} x_1' &= x_1 \ x_2' &= 3x_2. \end{aligned} \ ilde{f}(x') &= rac{1}{2} {x_1'}^2 + rac{9}{2} \left(rac{1}{3} {x_2'}
ight)^2 = rac{1}{2} {x_1'}^2 + rac{1}{2} {x_2'}^2 \,. \end{aligned}$$

Preconditioning



#### Preconditioning

- $ightharpoonup H_k$  symmetric positive definite.
- $ightharpoonup H_k = L_k L_k^T$
- Change of variables:

$$x' = L_k^T x$$

Steepest descent iteration:

$$x'_{k+1} = x'_k - \alpha_k \nabla \tilde{f}(x'_k)$$

$$\tilde{f}(x') = f(L_k^{-T} x') 
\nabla \tilde{f}(x') = L_k^{-1} \nabla f(L_k^{-T} x') 
x'_{k+1} = x'_k - \alpha_k L_k^{-1} \nabla f(L_k^{-T} x'_k) . 
L_k^T x_{k+1} = L_k^T x_k - \alpha_k L_k^{-1} \nabla f(x_k) 
x_{k+1} = x_k - \alpha_k L_k^{-T} L_k^{-1} \nabla f(x_k) 
x_{k+1} = x_k - \alpha_k H_k^{-1} \nabla f(x_k) .$$

#### Preconditioned steepest descent

$$x_{k+1} = x_k + \alpha_k d_k$$

with

$$d_k = -D_k \nabla f(x_k).$$

#### Preconditioner

Symmetric positive definite matrix. For instance,

$$D_k = H_k^{-1}.$$

#### Descent direction

$$\nabla f(x_k)^T d_k = -\nabla f(x_k)^T D_k \nabla f(x_k) < 0.$$

#### Inexact line search

#### Motivation

- ► Finding a local optimum along a direction takes a large effort, that may not be justified, as it has to be done at each iteration.
- What about choosing any step that would decrease the function?
- ▶ We first show that this does not always work.
- ▶ And then we propose something similar, that works.

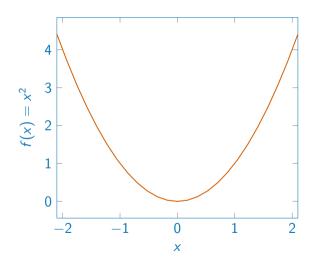
#### General idea

$$f: \mathbb{R}^n \to \mathbb{R}, \ x_k \in \mathbb{R}^n, \ d_k \in \mathbb{R}^n, \ \nabla f(x_k)^T d_k < 0.$$

$$x_{k+1} = x_k + \alpha_k d_k,$$

where  $\alpha_k$  is such that

$$f(x_{k+1}) < f(x_k).$$

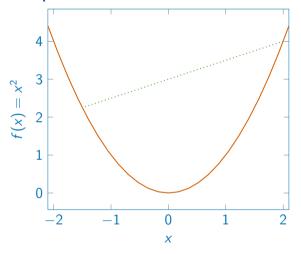


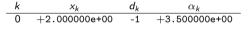
$$f(x) = x^{2}$$

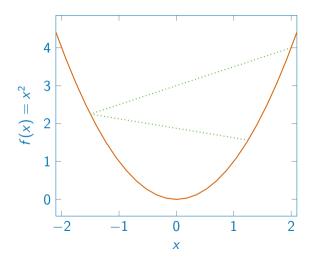
$$x_{0} = 2$$

$$d_{k} = -\operatorname{sgn}(x_{k})$$

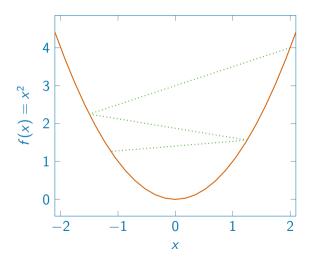
$$\alpha_{k} = 2 + 3(2^{-k-1}).$$



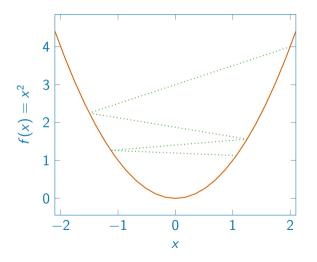




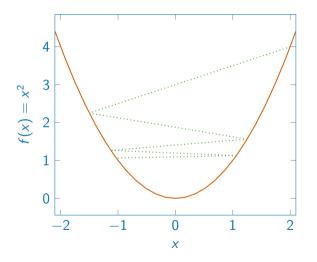
k	$\times_k$	$d_k$	$lpha_k$
0	+2.000000e+00	-1	+3.500000e+00
1	-1.500000e+00	1	+2.750000e+00



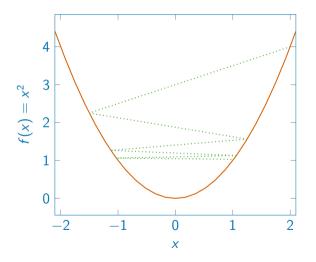
k	$\times_k$	$d_k$	$\alpha_k$
0	+2.000000e+00	-1	+3.500000e+00
1	-1.500000e+00	1	+2.750000e+00
2	$\pm 1.250000e \pm 00$	-1	$\pm 2.375000e \pm 00$



k	$\times_k$	$d_k$	$\alpha_{\pmb{k}}$
0	+2.000000e+00	-1	+3.500000e+00
1	-1.500000e+00	1	+2.750000e+00
2	+1.250000e+00	-1	+2.375000e+00
3	-1.125000e+00	1	+2.187500e+00

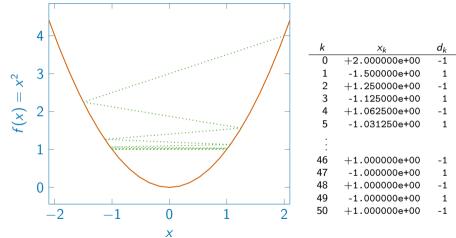


k	$\times_k$	$d_k$	$\alpha_k$
0	+2.000000e+00	-1	+3.500000e+00
1	-1.500000e+00	1	+2.750000e+00
2	+1.250000e+00	-1	+2.375000e+00
3	-1.125000e+00	1	+2.187500e+00
1	1 1 0625000+00	- 1	12 0027500+00



k	$\times_k$	$d_k$	$\alpha_k$
0	+2.000000e+00	-1	+3.500000e+00
1	-1.500000e+00	1	+2.750000e+00
2	+1.250000e+00	-1	+2.375000e+00
3	-1.125000e+00	1	+2.187500e+00
4	+1.062500e+00	-1	+2.093750e+00
5	-1.031250e+00	1	+2.046875e+00

# Example: one dimension



$\times_k$	$a_k$	$\alpha_k$
+2.000000e+00	-1	+3.500000e+00
-1.500000e+00	1	+2.750000e+00
+1.250000e+00	-1	+2.375000e+00
-1.125000e+00	1	+2.187500e+00
+1.062500e+00	-1	+2.093750e+00
-1.031250e+00	1	+2.046875e+00
+1.000000e+00	-1	+2.000000e+00
-1.000000e+00	1	+2.000000e+00
+1.000000e+00	-1	+2.000000e+00
-1.000000e+00	1	+2.000000e+00
+1.000000e+00	-1	+2.000000e+00
	+2.000000e+00 -1.500000e+00 +1.250000e+00 -1.125000e+00 +1.062500e+00 -1.031250e+00 +1.000000e+00 -1.000000e+00 -1.000000e+00 -1.000000e+00	+2.000000e+00 -1 -1.500000e+00 1 +1.250000e+00 -1 -1.125000e+00 -1 -1.031250e+00 1 +1.000000e+00 -1 -1.000000e+00 -1 -1.000000e+00 -1 -1.000000e+00 1

## Example: one dimension

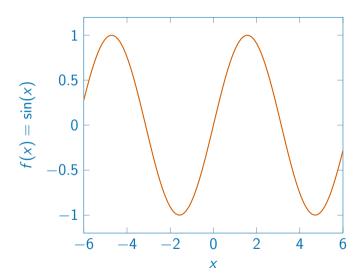
$$x_k = (-1)^k (1+2^{-k})$$
  $|x_{k+1}| < |x_k|$  therefore  $x_{k+1}^2 < x_k^2$  and  $f(x_{k+1}) < f(x_k)$ .  $\lim_{k \to \infty} x_k$  does not exists

Two subsequences converging to -1 and 1.

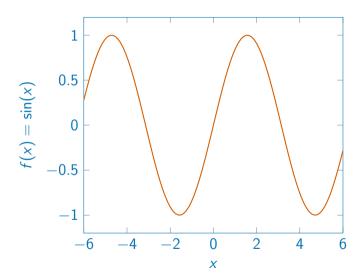
- Descent direction is a local concept.
- ► It assumes small steps.

- Descent direction is a local concept.
- lt assumes small steps.
- When steps are too long, Taylor's theorem does not apply anymore.

- Descent direction is a local concept.
- ► It assumes small steps.
- When steps are too long, Taylor's theorem does not apply anymore.
- Even ascent directions can decrease the function with long steps.



- Descent direction is a local concept.
- ► It assumes small steps.
- When steps are too long, Taylor's theorem does not apply anymore.
- Even ascent directions can decrease the function with long steps.
- ► Solution: avoid long steps.



# Second example

$$f(x) = x^{2}$$

$$x_{0} = 2$$

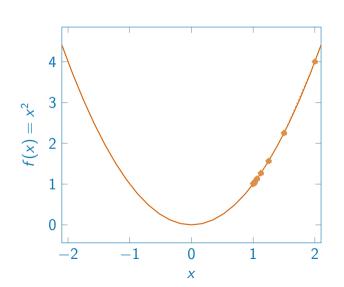
$$d_{k} = -1$$

$$\alpha_{k} = 2^{-k-1}$$

$$x_{k} = 1 + 2^{-k} > 0$$

$$x_{k+1} < x_{k} \iff f(x_{k+1}) < f(x_{k})$$

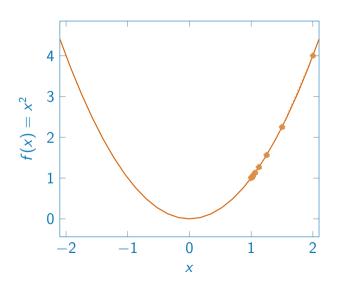
$$\lim_{k \to \infty} x_{k} = 1$$



Steps become smaller and smaller:

$$\lim_{k\to\infty}\alpha_k=\lim_{k\to\infty}2^{-k-1}=0.$$

- Although  $\alpha_k > 0$ , almost no progress can be made.
- ► Solution: avoid short steps.



### First Wolfe condition

### Motivation

- Our objective is to identify potential steps that can be done along a descent direction.
- ▶ We saw that the algorithm may fail to converge if the steps are too long.
- ▶ We provide here a characterization of the concept of "being too long".

### First idea

$$f: \mathbb{R}^n \to \mathbb{R}, \ x_k \in \mathbb{R}^n, \ d_k \in \mathbb{R}^n, \ \nabla f(x_k)^T d_k < 0.$$

### Require a decrease proportional to the step

$$f(x_k) - f(x_k + \alpha_k d_k) \ge \alpha_k \gamma$$
,

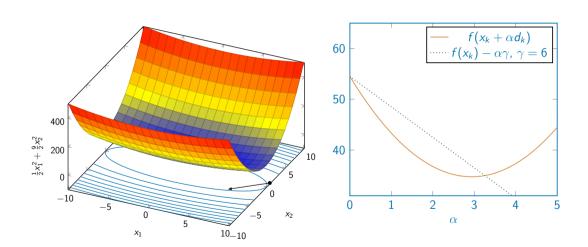
or

$$f(x_k + \alpha_k d_k) \leq f(x_k) - \alpha_k \gamma$$
,

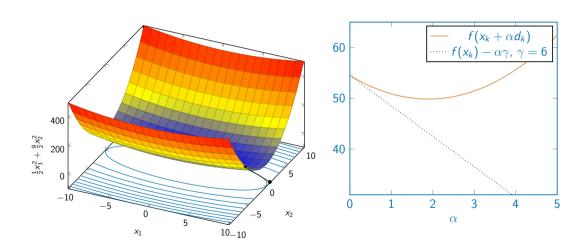
with

$$\gamma > 0$$
.

# Decrease proportional to the step



## But...



# Second idea





# More formally

$$f: \mathbb{R}^n \to \mathbb{R}, \ x_k \in \mathbb{R}^n, \ d_k \in \mathbb{R}^n, \ \nabla f(x_k)^T d_k < 0.$$

Choose  $\gamma$  according to the slope:

$$\gamma = -\beta_1 \nabla f(x_k)^T d_k, 0 < \beta_1 < 1$$

## First Wolfe condition

$$f(x_k + \alpha_k d_k) \le f(x_k) - \alpha_k \gamma,$$
  
 $\gamma = -\beta_1 \nabla f(x_k)^T d_k.$ 

$$f(x_k + \alpha_k d_k) \leq f(x_k) + \alpha_k \beta_1 \nabla f(x_k)^T d_k.$$

## Extreme cases are excluded

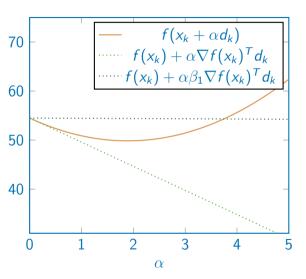
$$\beta_1 = 0$$

$$f(x_k + \alpha_k d_k) \leq f(x_k).$$

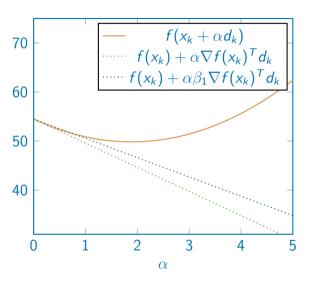
$$\beta_1 = 1$$

$$f(x_k + \alpha_k d_k) \leq f(x_k) + \alpha_k \nabla f(x_k)^T d_k$$
.

## $\beta_1 = 0.01$



## $\beta_1 = 0.8$



### Second Wolfe condition

#### Motivation

- Our objective is to identify valid steps along a descent direction.
- We have first identified a condition that prohibits steps that are too long.
- We are now deriving a condition that prohibits steps that are too short.

## Main idea

$$f: \mathbb{R}^n \to \mathbb{R}, \ x_k \in \mathbb{R}^n, \ d_k \in \mathbb{R}^n, \ \nabla f(x_k)^T d_k < 0.$$

### Main idea

$$f: \mathbb{R}^n \to \mathbb{R}, \ x_k \in \mathbb{R}^n, \ d_k \in \mathbb{R}^n, \ \nabla f(x_k)^T d_k < 0.$$

## Directional derivative along $x_k + \alpha d_k$

At  $\alpha = 0$ :

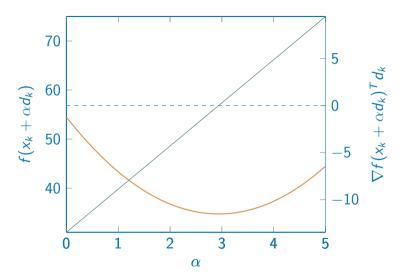
$$\nabla f(x_k)^T d_k < 0.$$

ightharpoonup At  $\alpha^*$ , first local minimum:

$$\nabla f(\mathbf{x}_k + \alpha^* d_k)^T d_k = 0.$$

▶ If we reach  $\alpha^*$ , the directional derivative increases by  $\nabla f(x_k)^T d_k$ .

## Directional derivative



### Second Wolfe condition

### Relative reduction of the derivative

$$\frac{\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k}{\nabla f(\mathbf{x}_k)^T \mathbf{d}_k} \le \beta_2$$

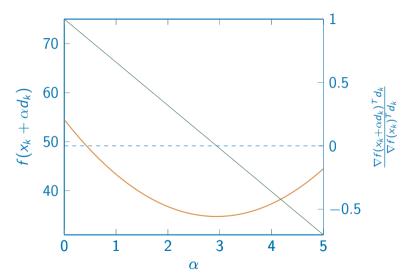
with  $0 < \beta_2 < 1$ .

## Sufficient progress

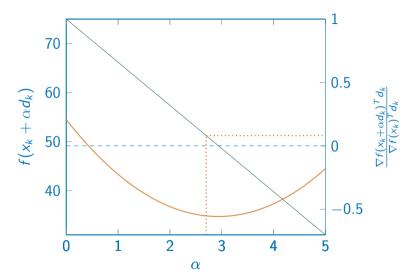
$$\nabla f(x_k + \alpha_k d_k)^T d_k \ge \beta_2 \nabla f(x_k)^T d_k$$

as  $\nabla f(x_k)^T d_k < 0$ .

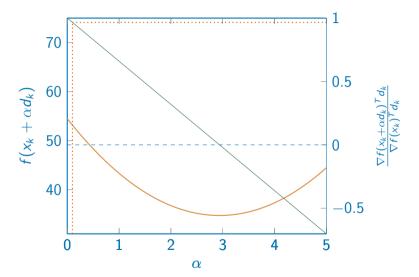
# Sufficient progress



# $\beta_2$ close to 0: larger steps



## $\beta_2$ close to 1: smaller steps



# Validity of the Wolfe conditions

#### Motivation

- First Wolfe condition forbids steps that are too long.
- Second Wolfe condition forbids steps that are too short.
- How do we guarantee that there exists steps that verify both?

### Wolfe conditions

### First Wolfe condition

$$f(x_k + \alpha_k d_k) \le f(x_k) + \alpha_k \beta_1 \nabla f(x_k)^T d_k,$$
  
$$0 < \beta_1 < 1.$$

### Wolfe conditions

### First Wolfe condition

$$f(x_k + \alpha_k d_k) \le f(x_k) + \alpha_k \beta_1 \nabla f(x_k)^T d_k,$$
  
$$0 < \beta_1 < 1.$$

### Second Wolfe condition

$$\frac{\nabla f(x_k + \alpha_k d_k)^T d_k}{\nabla f(x_k)^T d_k} \le \beta_2,$$
$$0 < \beta_2 < 1.$$

# Compatibility of the Wolfe conditions

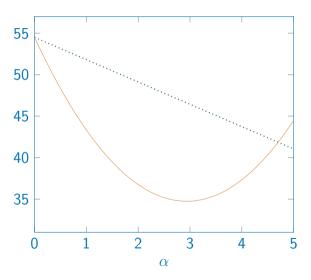
### Theorem 11.9

- ▶ If f is bounded from below along  $d_k$ ,
- if  $0 < \beta_1 < \beta_2 < 1$ ,
- there exists  $\alpha > 0$  verifying both conditions.

### First Wolfe condition

- ightharpoonup As  $d_k$  is a descent direction,
- ▶  $\exists \eta$  such that Wolfe 1 is verified for  $\alpha_k \leq \eta$ .
- ► See Theorem 2.11.

# Validity of the conditions



## Ideas of the proof

See Figure 11.15 p. 273

- ► Start with the line:  $f(x_k) + \alpha \beta_1 \nabla f(x_k)^T d_k$
- $\triangleright$   $\alpha_1$ : intersection between the line and the function.
- $ightharpoonup \alpha_1$  exists because the function is bounded from below.
- ▶ All  $\alpha \leq \alpha_1$  verify Wolfe 1.
- $\triangleright$   $\alpha_2$ : mean value theorem (tangent with same slope). Verifies Wolfe 1.
- ▶ As the slope is defined by  $\beta_1$ , we have

$$\beta_1 \nabla f(x_k)^T d_k = \nabla f(x_k + \alpha_2 d_k)^T d_k$$

so that

$$\beta_1 = \frac{\nabla f(x_k + \alpha_2 d_k)^T d_k}{\nabla f(x_k)^T d_k}$$

▶ If  $\beta_2 > \beta_1$ ,  $\alpha_2$  verifies Wolfe 2.

# Linesearch algorithm

### **Initialize:** $i \leftarrow 0$ , $\alpha_{\ell} \leftarrow 0$ , $\alpha_{r} \leftarrow +\infty$ .

### 1 repeat

5

10

11

if  $\alpha_i$  violates Wolfe 1 then the step is too long

if  $\alpha_i$  does not violate Wolfe 1 but violates Wolfe 2 then the step is too short

$$\alpha_r \leftarrow \alpha_i$$

$$\alpha_r \leftarrow \alpha_i, \\ \alpha_{i+1} \leftarrow \frac{\alpha_\ell + \alpha_r}{2}.$$

$$\alpha_{\ell} \leftarrow \alpha_{i}$$
.

if 
$$\alpha_r < +\infty$$
, then

$$-\infty$$
, ther

$$\alpha_{i+1} \leftarrow \frac{\alpha_{\ell} + \alpha_{r}}{2}$$
 else

## $\alpha_{i+1} \leftarrow 10\alpha_i$ .

$$i \leftarrow i + 1$$
.

12 until 
$$\alpha_i$$
 satisfies Wolfe 1 and 2

13  $\alpha^* \leftarrow \alpha$ :

# Finiteness of the line search algorithm

#### Motivation

- ► The linesearch algorithm is quite simple:
  - 1. Start with a candidate step.
  - 2. If it is too long, make it shorter.
  - 3. If it is too short, make it longer.
- ▶ We now need to verify that this process finishes in a finite number of steps.

### Wolfe conditions

### First Wolfe condition

$$f(x_k + \alpha_k d_k) \leq f(x_k) + \alpha_k \beta_1 \nabla f(x_k)^T d_k$$
.

### Second Wolfe condition

$$\nabla f(x_k + \alpha_k d_k)^T d_k \geq \beta_2 \nabla f(x_k)^T d_k.$$

### **Parameters**

$$0 < \beta_1 < \beta_2 < 1$$
.

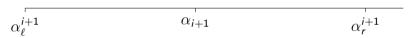
### Initialization



#### Initialization



### Wolfe 1 violated: step too long



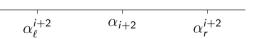
#### Initialization

$$\alpha_i^i$$

Wolfe 1 violated: step too long

$$\alpha_{\ell}^{i+1}$$
  $\alpha_{i+1}^{i+1}$   $\alpha_{r}^{i+1}$ 

Wolfe 2 violated: step too short



### Special case

- ▶ If the step is too short (Wolfe 2 violated), and
- $ightharpoonup \alpha_r = \infty$ ,
- ▶ then increase the step by an arbitrary factor  $\lambda > 1$ :

$$\alpha_{i+1} = \lambda \alpha_i.$$

## Properties

 $\alpha_{k}^{\prime}$ 

always verifies Wolfe 1: it is updated only when Wolfe 1 is verified and Wolfe 2 is violated.

## **Properties**

 $\alpha_{k}^{\prime}$ 

always verifies Wolfe 1: it is updated only when Wolfe 1 is verified and Wolfe 2 is violated.

 $\alpha_{\ell}^{i}$ 

always violates Wolfe 2 for the same reason.

## **Properties**

```
\alpha_{i}
```

always verifies Wolfe 1: it is updated only when Wolfe 1 is verified and Wolfe 2 is violated.

 $\alpha_\ell^i$ 

always violates Wolfe 2 for the same reason.

 $lpha_{\it r}^{\it i}$ 

always violates Wolfe 1 : starts from  $\infty$  and updated when Wolfe 1 is violated.

## Finiteness of $\alpha_i$

### Arguments

- ▶ Suppose, by contradiction, that  $\alpha_i \to \infty$ .
- ▶ It means that  $\alpha_r^i$  is always  $+\infty$ .
- As  $\alpha_r^i$  would be updated when Wolfe 1 is violated, it means that Wolfe 1 is never violated.
- Impossible, as the function is bounded from below.

## Finiteness of $\alpha_i$

### Arguments

- ▶ Suppose, by contradiction, that  $\alpha_i \to \infty$ .
- ▶ It means that  $\alpha_r^i$  is always  $+\infty$ .
- As  $\alpha_r^i$  would be updated when Wolfe 1 is violated, it means that Wolfe 1 is never violated.
- Impossible, as the function is bounded from below.

Consider only iterations when  $\alpha_r^i < \infty$ .

$$\alpha_{i+1} = \frac{\alpha_\ell^i + \alpha_r^i}{2}.$$

#### Theorem 11.10

- ▶ Suppose, by contradiction, that there is an infinite number of iterations.
- ► Therefore,

$$\lim_{i\to\infty}\alpha_r^i-\alpha_\ell^i=0.$$

Consequently,

$$\alpha^* = \lim_{i \to \infty} \alpha_r^i = \lim_{i \to \infty} \alpha_\ell^i = \lim_{i \to \infty} \alpha_i.$$

#### Theorem 11.10

▶ Wolfe 1 is verified for all  $\alpha_{\ell}^{i}$ . At the limit:

$$f(x_k + \alpha^* d_k) \leq f(x_k) + \alpha^* \beta_1 \nabla f(x_k)^T d_k.$$

▶ Wolfe 1 is violated by all  $\alpha_r^i$ . Therefore  $\alpha^* \neq \alpha_r^i$ .

$$f(x_k + \alpha_r^i d_k) > f(x_k) + \alpha_r^i \beta_1 \nabla f(x_k)^T d_k.$$

At the limit

$$f(x_k + \alpha^* d_k) \ge f(x_k) + \alpha^* \beta_1 \nabla f(x_k)^T d_k.$$

#### Theorem 11.10

$$f(\mathbf{x}_k + \alpha^* d_k) = f(\mathbf{x}_k) + \alpha^* \beta_1 \nabla f(\mathbf{x}_k)^T d_k.$$

▶ Wolfe 1 is violated by all  $\alpha_r^i$ . Therefore  $\alpha_r^i - \alpha^* > 0$ .

$$f(x_k + \alpha_r^i d_k) > f(x_k) + \alpha_r^i \beta_1 \nabla f(x_k)^T d_k.$$

$$f(x_k + \alpha_r^i d_k) > f(x_k + \alpha^* d_k) - \alpha^* \beta_1 \nabla f(x_k)^T d_k + \alpha_r^i \beta_1 \nabla f(x_k)^T d_k.$$

$$f(x_k + \alpha_r^i d_k) > f(x_k + \alpha^* d_k) + (\alpha_r^i - \alpha^*) \beta_1 \nabla f(x_k)^T d_k.$$

$$\frac{f(x_k + \alpha_r^i d_k) - f(x_k + \alpha^* d_k)}{\alpha_r^i - \alpha^*} > \beta_1 \nabla f(x_k)^T d_k.$$

$$\frac{f(x_k + \alpha_r^i d_k) - f(x_k + \alpha^* d_k)}{\alpha_r^i - \alpha^*} > \beta_1 \nabla f(x_k)^T d_k.$$

$$i \to \infty : \ \nabla f(x_k + \alpha^* d_k)^T d_k \ge \beta_1 \nabla f(x_k)^T d_k.$$
As  $\beta_2 > \beta_1$  and  $\nabla f(x_k)^T d_k < 0$ , we have
$$\nabla f(x_k + \alpha^* d_k)^T d_k > \beta_2 \nabla f(x_k)^T d_k.$$

#### Current result

$$\nabla f(\mathbf{x}_k + \alpha^* d_k)^T d_k > \beta_2 \nabla f(\mathbf{x}_k)^T d_k.$$

 $\alpha_{\ell}^{i}$  violates Wolfe 2

$$\nabla f(\mathbf{x}_k + \alpha_\ell^i d_k)^T d_k < \beta_2 \nabla f(\mathbf{x}_k)^T d_k, \ \forall i.$$

At the limit

$$\nabla f(x_k + \alpha^* d_k)^T d_k \leq \beta_2 \nabla f(x_k)^T d_k.$$

Contradiction

### Newton method with line search

#### Motivation

- ▶ We now go back to Newton's method.
- We have seen that, when it works, it converges fast.
- But it does not always work.
- Let's combine it with the descent methods framework to obtain an efficient algorithm.

## Comparison

#### Descent method

$$x_{k+1} = x_k - \alpha_k D_k \nabla f(x_k).$$

#### Newton's method

$$x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k).$$

## Comparison

#### Descent method

$$x_{k+1} = x_k - \alpha_k D_k \nabla f(x_k).$$

#### Newton's method

$$x_{k+1} = x_k - \nabla^2 f(x_k)^{-1} \nabla f(x_k).$$

### Newton's method is a descent method

- ightharpoonup if  $\alpha_k = 1$ ,
- ▶ if  $D_k = \nabla^2 f(x_k)^{-1}$  is positive definite.

### **Modifications**

## If $\nabla^2 f(x_k)^{-1}$ not positive definite

Choose another  $D_k$ :

- $\triangleright$   $D_k = I$ : bad idea. Slow convergence of steepest descent.
- $ightharpoonup D_k$  diagonal:

$$D_k(i,i) = \max\left(\varepsilon, \frac{\partial f^2}{\partial x_i^2}(x_k)\right)^{-1},$$

with  $\varepsilon > 0$ .

▶ Inflating  $\nabla^2 f(x_k)$ :

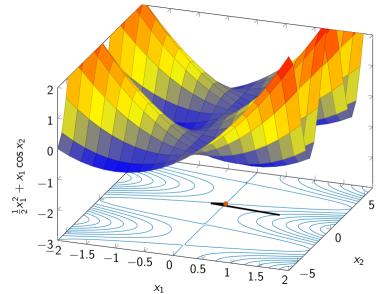
$$D_k = (\nabla^2 f(x_k) + \tau I)^{-1},$$

where  $\tau$  is calculated such that  $D_k$  is positive definite.

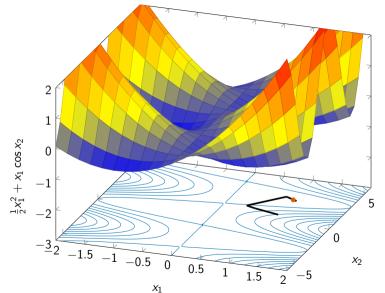
### Modifications

If  $\alpha_k = 1$  is not consistent with Wolfe conditions Apply the line search algorithm.

Example: Newton's method



# Example: Modified Newton's method



## **Iterations**

k	$f(x_k)$	$\ \nabla f(x_k)\ $	$\alpha_{\mathbf{k}}$	au
0	1.04030231e+00	1.75516512e+00		
1	2.34942031e-01	8.88574897e-01	1	1.64562250e+00
2	4.21849003e-02	4.80063696e-01	1	1.72091923e+00
3	-4.52738278e-01	2.67168927e-01	3	8.64490594e-01
4	-4.93913638e-01	1.14762780e-01	1	0.00000000e+00
5	-4.99982955e-01	5.85174623e-03	1	0.00000000e+00
6	-5.00000000e-01	1.94633135e-05	1	0.0000000e+00
7	-5.00000000e-01	2.18521663e-10	1	0.00000000e+00
8	-5.00000000e-01	1.22460635e-16	1	0.00000000e+00

# Fast and reliable



## Summary

- Optimality conditions: solve a system of nonlinear equations.
- Newton's method for optimization is fast, but unreliable.
- Other method: preconditioned steepest descent.
- Descent direction: preconditioned gradient.
- Step: Wolfe conditions.
- ▶ Newton method with line search.